

# Zero Shot Hashing

Shubham Pachori  
*Electrical Engineering*  
*Indian Institute of Technology Gandhinagar*  
*Gandhinagar, Gujarat 382355*  
*Email: shubham\_pachori@iitgn.ac.in*

Shanmuganathan Raman  
*Electrical Engineering &*  
*Computer Science and Engineering*  
*Indian Institute of Technology Gandhinagar*  
*Gandhinagar, Gujarat 382355*  
*Email: shanmuga@iitgn.ac.in*

## Abstract

*This paper provides a framework to hash images containing instances of unknown object classes. In many object recognition problems, we might have access to huge amount of data. It may so happen that even this huge data doesn't cover the objects belonging to classes that we see in our day to day life. Zero shot learning exploits auxiliary information (also called as signatures) in order to predict the labels corresponding to unknown classes. In this work, we attempt to generate the hash codes for images belonging to unseen classes, information of which is available only through the textual corpus. We formulate this as an unsupervised hashing formulation as the exact labels are not available for the instances of unseen classes. We show that the proposed solution is able to generate hash codes which can predict labels corresponding to unseen classes with appreciably good precision.*

## 1. Introduction

With billions of image-based data information added to social networking sites like Flickr, Instagram everyday, it has become challenge to accurately organize the data. One possible solution is to use hashing techniques with the smallest number of possible bits in order to reduce both storage requirement and query response time. The hashes corresponding to images thus obtained can be used for multiple computer vision tasks such as recognition, image retrieval and understanding. Hashing based approximate nearest neighbor(ANN) search methods have attracted a lot of attention in the past two decades. Apart from this, considerable amount of research has been done in the past few years on improving zero shot learning algorithms [ [24], [33], [10], [5], [17], [46]]. Zero shot learning requires one

to transfer knowledge learnt from the classes present in the training data to the classes which are not been observed yet. This knowledge is generally available in the form of signatures or attributes along with visual concept. Our main motivation for zero shot hashing (referred as ZSH in the rest of the paper) comes from the fact that humans learn to visualize an image from just the attributes present in it. For example, if we are given an dictionary where we store images, then we are likely to keep the image of liger in between the images of tiger and lion [2]. This has also been shown by Yosinski [44], where each filter present in CNN during training learns to detect different features in an image like clothes, face, numbers, etc. This cumulative knowledge helps to determine the classes to which the given image belongs to. The concept of attribute based learning has also been used in multiple instance learning, where bag of instances with same labels are created. These instances contain different attributes of an image such as strips, ears of leopard etc.

Features learnt from other modalities like text are transferred inductively to images using experience. The features extraction methods like CNN contain information about all the attributes present in a given image. However, they do not contain information about other modalities like text. Thus, we could keep an image of liger in between tiger and lion but we are allowed to create another class called "LIGER" in our dictionary. Moreover, it has been shown that using supervised information of an image along with its features could significantly improve the hash codes of image [43], [19], [23].

The primary contributions of this paper are listed below:

- 1) Incorporating zero shot learning framework in

induction based unsupervised hashing problem.

- 2) Learning correspondence between signatures of classes and images while jointly embedding them into a common space.
- 3) A novel but a simple approach to address out-of-sample extension problems associated with hashing images belonging to instances of unseen categories.

The outline of this paper is as follows. In section 2, we discuss the related works done previously and why our method is different from them. In section 3, we discuss the proposed methodology. In section 4, we evaluate our method and discuss the results obtained from our experiments. In section 5, we conclude our paper along with works in future.

## 2. Related Work

Many hashing methods have been proposed in the past, which can be categorized into two type - data dependent and data independent hashing. Locality sensitive hashing is the most popular data independent hashing technique, which uses randomized projections to generate hash functions and ensures high collision probability for similar data points. Variants of LSH [6] have been developed by taking different distance measures like Mahalanobis distance [13], kernel similarity ([12], [26]) and  $p$ -norm distance [1]. Other forms of LSH could be found in detail in [39]. In general data independent hashing techniques exploit long hashes and several hash tables to achieve better performance in terms of precision and recall, thus rendering them limited in use for large scale applications. On the contrary, data dependent techniques could be classified into two types - supervised and unsupervised hashing. These algorithms tend to exploit the available training data to generate short binary codes. Among the data dependent methods, PCA-based hashing ([8], [38]), supervised and semi-supervised hashing ([38], [19], [23], [11]) and graph based hashing ([20], [42]) techniques are quite popular.

It has been shown that leveraging non-linear manifold embedding techniques have helped in generating better pairwise affinity preserving dense binary codes. Among well-known hashing algorithms which utilize this idea is spectral hashing ([42], [41]), which uses the eigenfunctions of the Laplacian matrix to capture variation in the data. anchor graph hashing, popularly known as AGH ([20], [18]), uses anchor graphs for generating hash codes for training and out-of-sample data efficiently and effectively. Benefiting from the properties of manifold approaches, Inductive manifold hashing (IMH) has been proposed in [30],

which computes the manifold of a given data point according to the manifold of its neighbours. In [29], a solution has been proposed for preserving the inner-product similarities among raw vectors, while tackling maximum inner product search (MIPS) problem.

Apart from these, extensive research has been done in producing multimodal hash codes. Data available to us is in multiple information types and contains both text tags and visual concepts. The concept behind cross modal hashing (CMH) methods, in general, is to project the multimodal data in a common hamming space so that the distance between similar data in heterogeneous modalities are preserved. In [49], linear cross-modal hashing (LCMH) has been proposed, which tries to preserve inter-similarities between different modalities and intra-similarity within each modality. In [3], it was proposed to map the data from different modalities into a common subspace using projections learnt from collective matrix factorization techniques. Extending the technique of collective matrix factorization, latent semantic sparse hashing [47] exploits sparse coding to learn hash functions. In [37], semantic topic multimodal hashing (STMH) is proposed which generates each bit in hash code by finding whether a concept is available in the original data or not. For achieving this, it maps the learned multimodal semantic features into a common subspace by modeling text into multiple semantic concepts and corresponding images as latent semantic concepts. Supervised cross-modal hashing algorithms have also been proposed over time. In [14], spectral hashing has been incorporated with the multi-view case. In [45] semantic correlations are maximized and used to embed semantic labels into the training procedure. In [48], kernel-based supervised hashing for cross-view similarity search (KSH-CV) learns kernel hash functions using adaboost algorithm. To preserve the semantic similarity, [40] uses multi-class logistic regression to project heterogeneous data into a semantic space and uses a boosting framework to learn hash functions.

The difference between our approach and methods adopted in cross-modal hashing is that other methods assume that the information about a given class is present in both (textual and image) modalities while training the algorithm to generate hash codes for images. While in our case information in one of the mode (image) for unknown classes is absent during the training phase.

### 3. Proposed Method

Our methodology to produce hash codes for the images belonging to seen and unseen classes has been explained in this section. In section 3.1, we introduce notations that have been used throughout the paper. In section 3.2 and 3.3, we propose the approach to hash images belonging to seen classes. In section 3.4 and 3.5 we discuss the approach to hash images of unseen classes.

#### 3.1. Notations

In this paper, we denote the number of seen classes by  $n_s$  and the number of unseen classes by  $n_u$ . Let  $N_s$  and  $N_u$  denote the number of instances belonging to seen classes available to us during training and instances of unseen classes, respectively. Vector and its transpose are denoted by lower case bold Roman letters such as  $\mathbf{x}$  and  $\mathbf{x}^\top$  respectively. Uppercase bold roman letters, such as  $\mathbf{M}$ , denote matrices.  $\mathbb{1}$  represents the indicator function.  $\|\cdot\|_F^2$  represents the Frobenius norm.

#### 3.2. Creating anchor points

In the feature space, ideally the set of classes must form separate clusters such that the data points belonging to a certain class should belong to the cluster representing that class. Initially, the information of only seen classes are available to us. Thus our objective is to create clusters using these training instances which represent seen classes. This objective could be formulated in the same way as that of  $k$ -means clustering but with a little modification that we assign a penalty of  $\beta$  if the assigned cluster number for a particular image is different from its true label. The formulated objective function is shown in Eq.1:

$$\underset{\Pi, \mu_1, \dots, \mu_k}{\operatorname{argmin}} \sum_{n,k} \pi_{nk} \|\mathbf{x}_n - \mu_k\| + \beta \sum_{n=1}^{N_s} \mathbb{1}(\pi_n \neq \mathbf{y}_n) \quad (1)$$

where  $\mu_i$ 's are cluster centers and  $\Pi = [\pi_1, \dots, \pi_{N_s}]$  is cluster assignments in one-of- $K$  encoding format. This formulation is similar to the one given in [31]. Though, in their formulation, the authors had assumed that the number of unseen classes were available initially and the number of clusters they assigned were equal to  $n_s + n_u$ . We have not assumed that constraint and have chosen the number of clusters  $k$  to be equal to the number of unseen classes  $n_s$ .  $\pi_{nk}$  is equal to one, if the  $n$ th instance belongs to the  $k$ th cluster. In our experiments, we chose  $\beta = 0.9$ .

Exploiting EM algorithm,  $\mu_i$ 's and  $\Pi$  are updated iteratively and alternatively by optimizing the objective function. At each iteration,  $\mu_i$ 's are updated as given in Eq. 2

$$\mu_i = \frac{\sum_{n=1}^{n_s} \mathbb{1}(\pi_{ni} = 1) \mathbf{x}_n}{\sum_{n=1}^{n_s} \mathbb{1}(\pi_{ni} = 1)} \quad (2)$$

$\Pi$  is updated by assigning each instance to the cluster that minimizes the corresponding term.  $\mu_i$ 's are initialized as randomly chosen data points so that they are as far as possible from each other.

Let us call these cluster centers  $\mu_1, \mu_2, \dots, \mu_{n_s}$  as anchors in the rest of the paper. We also take the mean of the features of given instances corresponding to each of the seen classes and assign the anchor to the particular class according to the Euclidean distance with respect to the mean of features corresponding to the instances belonging to that class. We will embed these in the lower dimensional space using manifold learning. The number of dimensions in the lower dimensional manifold space is equal to the length of the hash code with which we want to hash an image.

#### 3.3. Producing hash codes for images belonging to the seen classes

Assuming that during training, instances correspond to only the images of seen classes, we use the cluster centers or anchors obtained by optimizing the Eq.1 to generate the hash codes for the images corresponding to the seen classes. Let us consider that we have a manifold-based low dimensional embedding  $\mathbf{M} := \{\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_{n_s}\}$  corresponding to the  $n_s$  cluster centers. Given features of an image  $\mathbf{x}_i$  belonging to the seen class, we aim to generate an embedding  $\mathbf{m}_i$  such that it preserves the local neighborhood relationship both in the feature and the embedded space with respect to anchors. To obtain the embedding  $\mathbf{m}_i$  in the manifold space, given the features  $\mathbf{x}_i$  of an image corresponding to seen classes, the objective function shown in Eq.3 is minimized.

$$O(\mathbf{m}_i) = \sum_{q=1}^{n_s} w(\mathbf{x}_q, \mathbf{x}_i) \|\mathbf{m}_i - \mathbf{m}_q\|^2 \quad (3)$$

where,  $w(\mathbf{x}_q, \mathbf{x}_i)$  captures the likelihood that the given image belongs to the  $q$ th cluster. This can be obtained by calculating the Euclidean distance of an image from anchors  $\mathbf{x}_q$  in the feature space. These distances from the cluster centers are then converted into probabilities

or weights to which cluster, a data point belongs to using the exponential function shown in Eq.4.

$$w(\mathbf{x}_q, \mathbf{x}_i) = \exp\left(\frac{-\|\mathbf{x}_i - \mathbf{x}_q\|^2}{\sigma^2}\right) \quad (4)$$

where,  $\mathbf{x}_i$  are features of  $i$ th image and  $\mathbf{x}_q$  is the  $q$ th anchor. Here,  $\sigma$  is a parameter. Differentiating  $O(\mathbf{m}_i)$  with respect to  $\mathbf{m}_i$  and equating it to zero, we obtain equation 6,

$$\left. \frac{\partial O(\mathbf{m}_i)}{\partial \mathbf{m}_i} \right|_{\mathbf{m}_i = \mathbf{m}_i^*} = 2 \sum_{q=1}^{n_s} w(\mathbf{x}_q, \mathbf{x}_i) (\mathbf{m}_i^* - \mathbf{m}_q) = 0, \quad (5)$$

$$\mathbf{m}_i^* = \frac{\sum_{q=1}^{n_s} w(\mathbf{x}_q, \mathbf{x}_i) \mathbf{m}_q}{\sum_{q=1}^{n_s} w(\mathbf{x}_q, \mathbf{x}_i)} \quad (6)$$

The proposed method here has been inspired from Shen et. al. [30], where they have provided an inductive formulation to obtain the embedding of any point using the linear combination of the base embeddings. We take the top  $s$  weights for our purpose and set the other weights to zero. We will call them as ‘nearest anchors’ for a given data point throughout the paper. Apart from this, we multiplied the  $i$ th weight by an exponential factor  $\omega^{-i}$ . For our experiments, we took  $\omega = 5$ . This is done so that in the manifold space, the distance between the given data point and the cluster assigned to it gets further decreased. That is the value of the weight with the highest value is further increased. Finally, these top  $s$  weights are re-normalized such that they sum to 1.

We use Eq.4 and Eq.6 to obtain the embedding of an image with features  $\mathbf{x}_i$  belonging to any of the seen classes. Finally, we obtain the hash codes for the given image by binarizing the embedding as shown in Eq.7.

$$h(\mathbf{x}_i) = \text{sign}(\mathbf{m}_i^*) \quad (7)$$

where  $\text{sign}(\cdot)$  is the element-wise sign function defined in Eq.8.

$$\text{sign}(k) = \begin{cases} 1, & \text{if } k \geq 0 \\ -1, & \text{otherwise} \end{cases} \quad (8)$$

### 3.4. Producing anchors for the unseen classes

Our main concept behind producing anchors for new unseen classes is that similar classes share similar attributes. For example, the classes ‘MONKEY’ and ‘CHIMPANZEE’ share many attributes in common

with each other, thus having high amount of similarity with each other. This could be inferred from 1. Thus, labels for unseen categories could be embedded using the description about how similar they are to the seen classes [46]. We utilize the information of similarity between attributes of the classes to obtain the embeddings of the unseen classes. Thus, whenever we learn about the information (in terms of attributes) about new class by any means like textual description, we create an anchor corresponding to it. To obtain the anchor we use the cosine similarity measure between the attributes of a new class with respect to the classes, anchor of which has been obtained.

$$w(C_i, C_j) = \frac{\mathbf{a}_i \cdot \mathbf{a}_j}{\|\mathbf{a}_i\| \|\mathbf{a}_j\|} \quad (9)$$

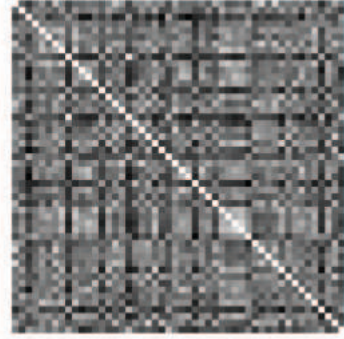


Figure 1: Cosine similarity between different classes of AWA dataset.

Here  $C_i$  and  $C_j$  are two classes between which cosine similarity is calculated and  $\mathbf{a}_i$  and  $\mathbf{a}_j$  are their corresponding attributes in the binary form.

We then use the Eq.6 to inductively obtain the embedding of the given class in the manifold space or anchor as shown in Eq.10.

$$\mathbf{m}_{C_i} = \frac{\sum_{q=1}^{n_s} w(C_q, C_i) \mathbf{m}_q}{\sum_{q=1}^{n_s} w(C_q, C_i)} \quad (10)$$

where,  $\mathbf{m}_{C_i}$  is the embedding in the manifold space for the unseen class  $C_i$ . This anchor is then added to our set of base anchors.

### 3.5. Generating hash codes for images of unseen classes

To produce the hash codes for images of unseen classes, we must embed the semantic information of

different classes i.e., attributes in a common space. The framework proposed by [27] is computationally cheaper, simple and provides a closed form solution of the problem. These are the main reasons due to which we adopt their approach.

**3.5.1. Zero-shot learning for images of unseen classes.** Let us assume that for each of the  $n_s$  classes at training stage, we have a signature vector of size  $\mathbf{a}$  such that each element of  $\mathbf{a}$  lies in  $[0, 1]$ . Signatures are represented in a matrix form as  $\mathbf{S} \in [0, 1]^{a \times n_s}$ . Let us denote all the instances available at training stage by a matrix  $\mathbf{X} \in \mathbb{R}^{N_s \times d}$ , where  $d$  is the length of feature vector of each instance. All instances are labeled in one hot encoding format, with ground truth labels of each of these instances are represented using as  $\mathbf{Y} \in \{-1, 1\}^{N_s \times n_s}$ , with positive entry indicating the class to which the instance belongs to. To learn a predictor corresponding to the  $n_s$  training classes, the following objective function is optimized:

$$\underset{\mathbf{V} \in \mathbb{R}^{d \times a}}{\operatorname{argmin}} L(\mathbf{XVS}, \mathbf{Y}) + \Omega(\mathbf{V}) \quad (11)$$

Here  $\mathbf{V} \in \mathbb{R}^{d \times a}$  embeds the semantic information in the form of attributes with image features. The regularizer chosen is of the following form shown in the equation below.

$$\Omega(\mathbf{V}; \mathbf{S}, \mathbf{X}) = \gamma \|\mathbf{VS}\|_F^2 + \lambda \|\mathbf{XV}\|_F^2 + \alpha \|\mathbf{V}\|_F^2 \quad (12)$$

The first term of the regularizer controls attribute signature so that their representations have a similar Euclidean norm on the feature space. The second term of the regularizer checks that the approach is invariant enough to be generalized to other test feature distribution by bounding the variance of representation of instances on attribute space. The third term penalises the Frobenius norm of the weight matrix to be learned.

The scalars  $\gamma, \lambda$  and  $\alpha$  are the hyper-parameters. If following choices are made

- 1)  $L(\mathbf{XVS}, \mathbf{Y}) = \|\mathbf{XVS} - \mathbf{Y}\|_F^2$
- 2)  $\alpha = \gamma\lambda$

then a closed solution of Eq.11 could be obtained as follows:

$$\mathbf{V} = (\mathbf{X}^\top \mathbf{X} + \gamma \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{YS}^\top (\mathbf{SS}^\top + \lambda \mathbf{I})^{-1} \quad (13)$$

At the testing stage, we are provided with signatures  $\mathbf{S}' \in [0, 1]^{a \times n_u}$  of unseen classes  $n_u$ . The probability that new instance from unseen class with feature vector  $\mathbf{x}_u$  belongs to the  $i$ th class  $C_i$  is calculated as:

$$w(\mathbf{x}_u, C_i) = \mathbf{x}_u \mathbf{VS}'_i \quad (14)$$

**3.5.2. Generating hash codes for images of unseen classes.** Once the probability to which class an image belongs to is calculated, we then use Eq.14 and Eq.6 to inductively produce hash codes for the images belonging to the unseen classes. Here also, we take the  $s$  nearest anchors for our purpose and multiplied the  $i$ th weight by an exponential factor  $\omega^{-i}$  before renormalizing these top  $s$  weights. The manifold embedding  $\mathbf{m}_{\mathbf{x}_u}$  of any instance with features  $\mathbf{x}_u$  from any of the unseen class is thus calculated as:

$$\mathbf{m}_{\mathbf{x}_u}^* = \frac{\sum_{i=1}^{n_u} w(\mathbf{x}_u, C_i) \mathbf{m}_{C_i}}{\sum_{i=1}^{n_u} w(\mathbf{x}_u, C_i)} \quad (15)$$

where  $\mathbf{m}_{C_i}$  is the manifold embedding of the cluster center of class  $C_i$ . To create the hash code for the given instance we binarize the hash code by taking using the  $\operatorname{sign}(\cdot)$  function.

## 4. Experimental Results

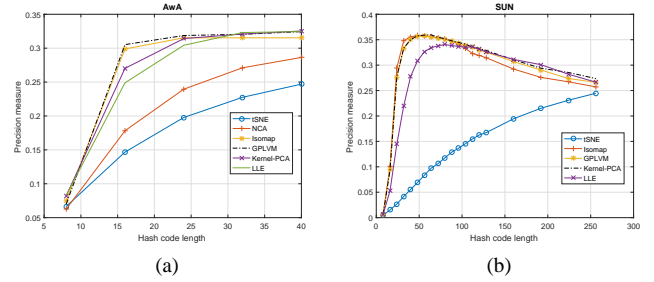


Figure 2: Comparison of different methods on AwA (left) and SUN (right) datasets based on precision measure for varying code lengths with hamming radius 2.

We performed our experiments on two datasets: the SUN scene attributes dataset [25] and the Animals with Attributes dataset (AwA) [15]. AwA dataset consists of 30K instances. Images are categorized into one of total 50 different classes each with binary attributes of 85 dimensions. For AwA dataset, attributes are provided for each class in the dataset. We used DeCAF [4] features for AwA dataset. SUN dataset consists of 14,340 images each categorized into one of total 717 different classes with 20 samples for each class. Each instance has attributes of 102 dimension with value in  $[0, 1]$ . For SUN dataset, attribute signature of each class is calculated by averaging the attribute signature of the instances belonging to that class. For SUN dataset, we

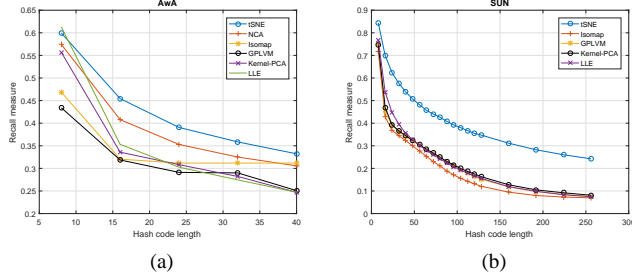


Figure 3: Comparison of different methods on AWA (left) and SUN (right) datasets based on recall measure for varying code lengths with hamming radius 2.

obtained deep features using VGG with 19-layer network [32] using MatConvNet [36]. We used six types of dimensionality reduction techniques in our work: Local Linear Embedding (LLE) [28], t-Distributed Stochastic Neighbor Embedding (t-SNE) [21], Isomap [34], Kernel-PCA [9], Gaussian Process latent variable model (GPLVM) [16] and Neighbourhood components analysis (NCA) [7]. We used the matlab toolbox for dimensionality reduction [35] to embed anchors in manifolds. We performed two sets of experiments. For each experiment, we ran 30 trials and present the averaged results. In one set of experiments, we found the accuracy by finding the Hamming distance of the hash codes with respect to the hash code of the anchors which are generated by binarizing their embeddings. The instance was assigned to that anchor for which the Hamming distance of the instance with respect to that anchor is the lowest. We exploit the label of each image for the ground truth. In the second set of experiments, we measure the performance of our algorithm using mean of average precision (MAP), precision and recall curves. We also show the hash lookup results using F1 measure [22] which is given as follows:

$$F1 = 2 \times \frac{(\text{precision} \cdot \text{recall})}{\text{precision} + \text{recall}}$$

For evaluating F1 measure, we used Hamming distance of 2 units throughout in the above experiments. For evaluating the precision, recall and F1 measures, we divided the dataset into two sets of quarter and three-quarter size, after obtaining the hash codes of each instance in our dataset. Then for each sample in the smaller set, we find all the samples in the larger set which are at hamming distance of 2 or less than 2 units from it. We exploit the ground truths provided to us for the images in the dataset to obtain precision, recall, F1-measure and MAP measure. We used  $s = 5$  in all the experiments unless specified explicitly. We

also measure the time computed by each method for embedding anchors (Table.1), producing hash code for each instance in the training set (Table.2) and test set (Table.3).

**Results on AWA dataset:** We randomly split the dataset into training and testing part with 40 classes in our training set and rest 10 classes for testing classes. We determined the hyper parameters by randomly taking 20 percent of the training dataset (8 classes) as our validation dataset, which were later combined with training set after fine tuning the model. For AWA dataset, the hyperparameters obtained are  $\gamma = 10$  and  $\lambda = 100$ . We tested our hash codes by taking hashes of bit size 8, 16, 24, 32, and 40 for AWA dataset. The maximum hash code length is equal to the number of seen classes as manifold embedding is done using the anchors which are equal to the number of seen classes. This is because, t-SNE and other manifold based techniques do not operate if the number of dimensions is less than number of data points as they are dependent on PCA framework.

**Results on SUN dataset:** We randomly split the dataset into training and testing part with 667 classes in our training set and rest 50 classes for testing classes. We determined the hyper parameters by randomly taking 20 percent of the training dataset (134 classes) as our validation dataset, which were later combined with training set after fine tuning the model. For the SUN dataset, the hyperparameters obtained are  $\gamma = 0.01$  and  $\lambda = 1$ . We tested our hash codes by taking hashes of bit size 8, 16, 24, 32, 40, 48, 56, 64, 72, 80, 88, 96, 104, 112, 120, 128, 160, 192, 224, and 256 for the SUN dataset.

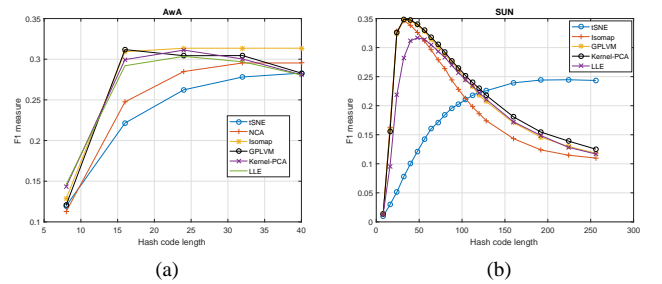


Figure 4: Comparison of different methods on AWA (left) and SUN (right) datasets based on F1 measure for varying code lengths with hamming radius 2.

From Fig. 5, it can be observed that for AWA dataset, with increasing hash code length, MAP measure increases for all the embeddings. But for the SUN dataset, as the hash code length increases, MAP measure for LLE based hashing decreases. For both



Table 1: Time (in sec) taken by each method to embed anchors in low dimensional space.

Method	AwA Dataset			SUN Dataset		
	24-bit	32-bit	40-bit	32-bit	64-bit	128-bit
tSNE-ZSH	0.117	0.128	0.128	6.893	7.198	8.393
NCA-ZSH	5.809	7.088	3.999	526.707	532.130	518.442
Isomap-ZSH	0.041	0.038	0.034	1.307	1.157	1.138
GPLVM-ZSH	0.027	<b>0.011</b>	0.083	<b>0.906</b>	<b>0.881</b>	<b>0.869</b>
Kernel PCA-ZSH	<b>0.022</b>	0.019	<b>0.012</b>	1.474	2.127	1.628
LLE-ZSH	16.272	12.367	13.193	15.213	14.313	14.647

Table 2: Time (in sec) taken by each method to produce hash codes for images in the training dataset.

Method	AwA Dataset			SUN Dataset		
	24-bit	32-bit	40-bit	32-bit	64-bit	128-bit
tSNE-ZSH	$1.271 \times 10^{-5}$	$1.213 \times 10^{-5}$	$1.197 \times 10^{-5}$	<b><math>2.157 \times 10^{-5}</math></b>	<b><math>2.184 \times 10^{-5}</math></b>	<b><math>2.129 \times 10^{-5}</math></b>
NCA-ZSH	$1.056 \times 10^{-5}$	$1.193 \times 10^{-5}$	$1.148 \times 10^{-5}$	$3.68 \times 10^{-5}$	$3.536 \times 10^{-5}$	$3.671 \times 10^{-5}$
Isomap-ZSH	<b><math>9.062 \times 10^{-6}</math></b>	<b><math>9.651 \times 10^{-6}</math></b>	$1.102 \times 10^{-5}$	$8.476 \times 10^{-5}$	$8.268 \times 10^{-5}$	$8.360 \times 10^{-5}$
GPLVM-ZSH	$1.334 \times 10^{-5}$	$1.143 \times 10^{-5}$	$1.178 \times 10^{-5}$	$9.053 \times 10^{-5}$	$9.297 \times 10^{-5}$	$9.104 \times 10^{-5}$
Kernel PCA-ZSH	$1.099 \times 10^{-5}$	$1.466 \times 10^{-5}$	<b><math>1.078 \times 10^{-5}</math></b>	$8.452 \times 10^{-5}$	$8.532 \times 10^{-5}$	$8.051 \times 10^{-5}$
LLE-ZSH	$1.472 \times 10^{-5}$	$1.152 \times 10^{-5}$	$1.288 \times 10^{-5}$	$3.299 \times 10^{-5}$	$3.138 \times 10^{-5}$	$3.287 \times 10^{-5}$

Table 3: Time (in sec) taken by each method to produce hash codes for images in the testing dataset.

Method	AwA Dataset			SUN Dataset		
	24-bit	32-bit	40-bit	32-bit	64-bit	128-bit
tSNE-ZSH	<b><math>2.108 \times 10^{-4}</math></b>	$2.118 \times 10^{-4}$	$2.122 \times 10^{-4}$	$7.017 \times 10^{-4}$	$7.031 \times 10^{-4}$	$7.019 \times 10^{-4}$
NCA-ZSH	$2.196 \times 10^{-4}$	$2.149 \times 10^{-4}$	$2.223 \times 10^{-4}$	$7.918 \times 10^{-4}$	$7.682 \times 10^{-4}$	$7.533 \times 10^{-4}$
Isomap-ZSH	$3.299 \times 10^{-4}$	$3.209 \times 10^{-4}$	$3.255 \times 10^{-4}$	<b><math>6.867 \times 10^{-4}</math></b>	<b><math>6.857 \times 10^{-4}</math></b>	<b><math>6.832 \times 10^{-4}</math></b>
GPLVM-ZSH	$3.483 \times 10^{-4}$	$3.418 \times 10^{-4}$	$3.380 \times 10^{-4}$	$8.368 \times 10^{-4}$	$8.431 \times 10^{-4}$	$9.284 \times 10^{-4}$
Kernel-PCA-ZSH	$4.178 \times 10^{-4}$	$4.031 \times 10^{-4}$	$3.696 \times 10^{-4}$	$8.443 \times 10^{-4}$	$8.273 \times 10^{-4}$	$8.760 \times 10^{-4}$
LLE-ZSH	$2.124 \times 10^{-4}$	<b><math>2.092 \times 10^{-4}</math></b>	<b><math>2.106 \times 10^{-4}</math></b>	$8.698 \times 10^{-4}$	$8.554 \times 10^{-4}$	$8.756 \times 10^{-4}$

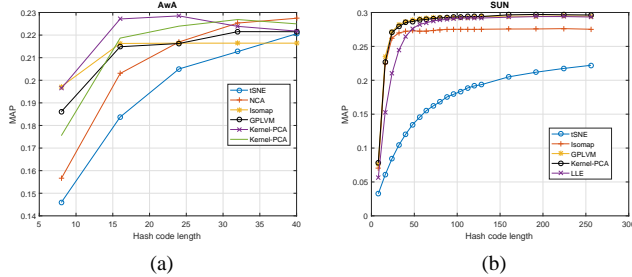


Figure 5: Comparison of different methods on AwA (left) and SUN (right) datasets based on MAP for varying code lengths with hamming radius 2.

the datasets, we observe that Kernel-PCA and NCA embedding based hashing methods perform better than the other methods. t-SNE embedding based hashing method gives poorer results when the hash code length is small but its performance significantly improves as the code length increases. From Fig.2, it can be observed that Kernel PCA and GPLVM based embedding methods perform superior than the other techniques for both datasets in term of precision for the hash codes of small length. From the precision measure of the SUN dataset, it can be observed that precision of methods

except t-SNE decreases after the hash code length of 56 but it increases for t-SNE continuously with increasing code length for both the datasets.

We also observe from the Fig. 6, that for AwA dataset, the accuracy for training data decreases with increasing the hash code length. While for the SUN dataset, the training data initially decreases but increases again and then saturates. For t-SNE based embedding, we observe that the training data accuracy decreases continuously. A possible reason could be that with increasing the dimension of t-SNE, it becomes less discrete and thus the hash code it generates becomes less distinct with increase in the length of hash code and hence many instances share similar hash codes. That is why, its recall is higher than other methods as it can be seen from Fig. 3 while its precision rate being low (Fig. 2). Similar reasons could be given regarding the performance of our method on accuracy of instances belonging to testing classes. The reason for this dramatic decrease in the performance of hash look up in Fig.4, Fig.6 and Fig.7 is that hamming spaces become sparser as we increase the hash code length.

We also evaluate our method by varying the number of nearest anchors used to obtain the embedding of data points and plot the MAP (Fig.9) and F1 measure

(Fig.8) curves. For this, we kept the hash code length fixed to 32 bits. We can notice from the plots that the performance of the proposed methods (except t-SNE) do not change significantly by varying the number of nearest anchors for any of the manifold embedding significantly for any of the two datasets. We see that LLE based embedding ZSH algorithm performs significantly better than its counterparts for both the datasets. While t-SNE based embedding hashing technique performs poorly in terms of both F1 and MAP measure. For t-SNE embedding based hashing, its performance increases continuously as we increase the number of nearest anchors.

## 5. Conclusions

We have proposed a hashing algorithm in the zero shot learning framework. Once the manifold embeddings are obtained corresponding to the training classes, our hashing formulation requires linear time for hashing all the training instances ( $O(n)$ ). We used different non-parametric dimensionality reduction techniques to preserve the data distribution in the original feature space. The proposed framework exploits the information of similarity between classes to inductively generate hash codes of images belonging to the unseen classes. One advantage of this type of hashing is that if an image (for e.g. images of claws of eagle) belonging to a seen class (eagle) shares a high similarity with an unseen class (hawk), its hash code will also be similar to the instances belonging to the unseen classes and we could still retrieve that image while querying for the unseen class. This is due to the fact that in our method the anchor corresponding to the unseen class is close to the anchor corresponding to the seen class. We also provided the methodology to generate hash codes of out-of-sample data.

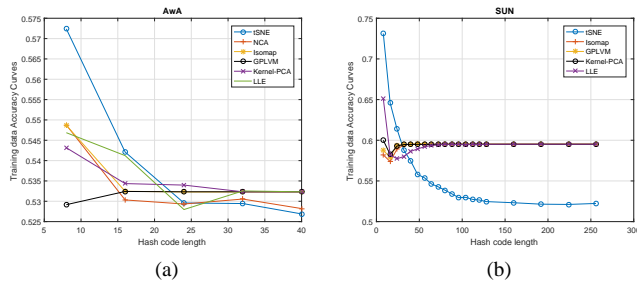


Figure 6: Comparison of different methods on Awa (left) and SUN (right) datasets based on training data accuracy for varying code lengths with hamming radius 2.

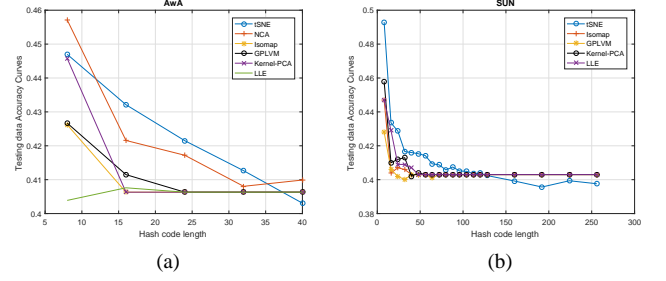


Figure 7: Comparison of different methods on Awa (left) and SUN (right) datasets based on testing data accuracy for varying code lengths with hamming radius 2.

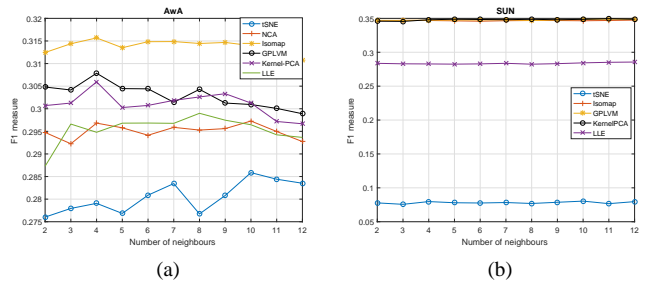


Figure 8: Comparison of different methods on Awa (left) and SUN (right) datasets based on F1 measure for varying number of nearest anchors.

## 6. Future Work

In the proposed method, we have not used significant amount of non-linearity for ranking the image features to the classes to which they belong to. In future, we could leverage deep learning methods which have achieved recently huge success in both fields of hashing and non-linear dimensionality reduction. Apart

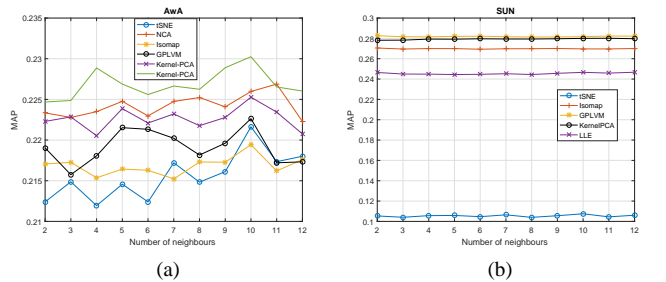


Figure 9: Comparison of different methods on Awa (left) and SUN (right) datasets based on MAP for varying number of nearest anchors.



from this, deep learning techniques have achieved positive results in embedding different modalities to a common space. Moreover, recently in [17], authors have utilized deep learning framework for learning joint latent space. This work has inspired us to use deep networks for improving the zero shot hashing framework in the future.

## Acknowledgments

The authors would like to thank Rajendra Nagar and Aalok Gangopadhyay for helpful discussions.

## References

- [1] M. Datar, N. Immorlica, P. Indyk, and V. S. Mirrokni. Locality-sensitive hashing scheme based on p-stable distributions. In *Proceedings of the twentieth annual symposium on Computational geometry*, pages 253–262. ACM, 2004.
- [2] R. Dawkins and D. McKean. *The Illustrated Magic of Reality: How We Know What’s Really True*. Simon and Schuster, 2012.
- [3] G. Ding, Y. Guo, and J. Zhou. Collective matrix factorization hashing for multimodal data. In *CVPR*. IEEE, 2014.
- [4] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition.
- [5] M. Elhoseiny, B. Saleh, and A. Elgammal. Write a classifier: Zero-shot learning using purely textual descriptions. In *ICCV*. IEEE, 2013.
- [6] A. Gionis, P. Indyk, R. Motwani, et al. Similarity search in high dimensions via hashing. In *VLDB*, volume 99, pages 518–529, 1999.
- [7] J. Goldberger, G. E. Hinton, S. T. Roweis, and R. Salakhutdinov. Neighbourhood components analysis. In *NIPS*, 2004.
- [8] Y. Gong, S. Lazebnik, A. Gordo, and F. Perronnin. Iterative quantization: A procrustean approach to learning binary codes for large-scale image retrieval. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(12):2916–2929, 2013.
- [9] J. Ham, D. D. Lee, S. Mika, and B. Schölkopf. A kernel view of the dimensionality reduction of manifolds. In *ICML*. ACM, 2004.
- [10] E. Kodirov, T. Xiang, Z. Fu, and S. Gong. Unsupervised domain adaptation for zero-shot learning. In *ICCV*. IEEE, 2015.
- [11] B. Kulis and T. Darrell. Learning to hash with binary reconstructive embeddings. In *NIPS*, 2009.
- [12] B. Kulis and K. Grauman. Kernelized locality-sensitive hashing for scalable image search. In *ICCV*. IEEE, 2009.
- [13] B. Kulis, P. Jain, and K. Grauman. Fast similarity search for learned metrics. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(12):2143–2157, 2009.
- [14] S. Kumar and R. Udupa. Learning hash functions for cross-view similarity search. In *IJCAI*, 2011.
- [15] C. H. Lampert, H. Nickisch, and S. Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *CVPR*. IEEE, 2009.
- [16] N. D. Lawrence. Gaussian process latent variable models for visualisation of high dimensional data. 2004.
- [17] J. Lei Ba, K. Swersky, S. Fidler, et al. Predicting deep zero-shot convolutional neural networks using textual descriptions. In *ICCV*. IEEE, 2015.
- [18] W. Liu, C. Mu, S. Kumar, and S.-F. Chang. Discrete graph hashing. In *NIPS*, 2014.
- [19] W. Liu, J. Wang, R. Ji, Y.-G. Jiang, and S.-F. Chang. Supervised hashing with kernels. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 2074–2081. IEEE, 2012.
- [20] W. Liu, J. Wang, S. Kumar, and S.-F. Chang. Hashing with graphs. In *ICML*, 2011.
- [21] L. v. d. Maaten and G. Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008.
- [22] C. D. Manning, P. Raghavan, and H. Schütze. Cambridge university press; 2008. *Introduction to Information Retrieval*, pages 158–163.
- [23] M. Norouzi and D. M. Blei. Minimal loss hashing for compact binary codes. In *ICML*, 2011.
- [24] M. Palatucci, D. Pomerleau, G. E. Hinton, and T. M. Mitchell. Zero-shot learning with semantic output codes. In *NIPS*, 2009.
- [25] G. Patterson and J. Hays. Sun attribute database: Discovering, annotating, and recognizing scene attributes. In *CVPR*. IEEE, 2012.
- [26] M. Raginsky and S. Lazebnik. Locality-sensitive binary codes from shift-invariant kernels. In *NIPS*, 2009.
- [27] B. Romera-Paredes and P. Torr. An embarrassingly simple approach to zero-shot learning. In *ICML*, 2015.
- [28] S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000.

- [29] F. Shen, W. Liu, S. Zhang, Y. Yang, and H. T. Shen. Learning binary codes for maximum inner product search. In *ICCV*. IEEE, 2015.
- [30] F. Shen, C. Shen, Q. Shi, A. Van Den Hengel, and Z. Tang. Inductive hashing on manifolds. In *CVPR*, 2013.
- [31] S. M. Shojaei and M. S. Baghshah. Semi-supervised zero-shot learning by a clustering-based approach. *arXiv preprint arXiv:1605.09016*, 2016.
- [32] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [33] R. Socher, M. Ganjoo, C. D. Manning, and A. Ng. Zero-shot learning through cross-modal transfer. In *NIPS*, 2013.
- [34] J. B. Tenenbaum, V. De Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
- [35] L. Van Der Maaten, E. Postma, and J. Van den Herik. Dimensionality reduction: a comparative. *J Mach Learn Res*, 10:66–71, 2009.
- [36] A. Vedaldi and K. Lenc. Matconvnet: Convolutional neural networks for matlab. In *Proceedings of the 23rd ACM international conference on Multimedia*, pages 689–692. ACM, 2015.
- [37] D. Wang, X. Gao, X. Wang, and L. He. Semantic topic multimodal hashing for cross-media retrieval. In *IJCAI*, 2015.
- [38] J. Wang, S. Kumar, and S.-F. Chang. Semi-supervised hashing for large-scale search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(12):2393–2406, 2012.
- [39] J. Wang, H. T. Shen, J. Song, and J. Ji. Hashing for similarity search: A survey. *arXiv preprint arXiv:1408.2927*, 2014.
- [40] K. Wang, J. Tang, N. Wang, and L. Shao. Semantic boosting cross-modal hashing for efficient multimedia retrieval. *Information Sciences*, 330:199–210, 2016.
- [41] Y. Weiss, R. Fergus, and A. Torralba. Multidimensional spectral hashing. In *ECCV*. Springer, 2012.
- [42] Y. Weiss, A. Torralba, and R. Fergus. Spectral hashing. In *NIPS*, 2009.
- [43] R. Xia, Y. Pan, H. Lai, C. Liu, and S. Yan. Supervised hashing for image retrieval via image representation learning. In *AAAI*, volume 1, page 2, 2014.
- [44] J. Yosinski, J. Clune, A. Nguyen, T. Fuchs, and H. Lipson. Understanding neural networks through deep visualization. *arXiv preprint arXiv:1506.06579*, 2015.
- [45] D. Zhang and W.-J. Li. Large-scale supervised multi-modal hashing with semantic correlation maximization. In *AAAI*, volume 1, page 7, 2014.
- [46] Z. Zhang and V. Saligrama. Zero-shot learning via semantic similarity embedding. In *ICCV*. IEEE, 2015.
- [47] J. Zhou, G. Ding, and Y. Guo. Latent semantic sparse hashing for cross-modal similarity search. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*, pages 415–424. ACM, 2014.
- [48] J. Zhou, G. Ding, Y. Guo, Q. Liu, and X. Dong. Kernel-based supervised hashing for cross-view similarity search. In *ICME*. IEEE, 2014.
- [49] X. Zhu, Z. Huang, H. T. Shen, and X. Zhao. Linear cross-modal hashing for efficient multimedia search. In *Proceedings of the 21st ACM international conference on Multimedia*, pages 143–152. ACM, 2013.